

Parallel Programming in CUDA C

Instructor

Dr. B. Krishna Priya

Outline

- CUDA parallel programming
- Summing vectors
- A fun example

CUDA Parallel Programming

Summing Vectors:

- Imagine having two lists of numbers where we want to sum corresponding elements of each list and store the

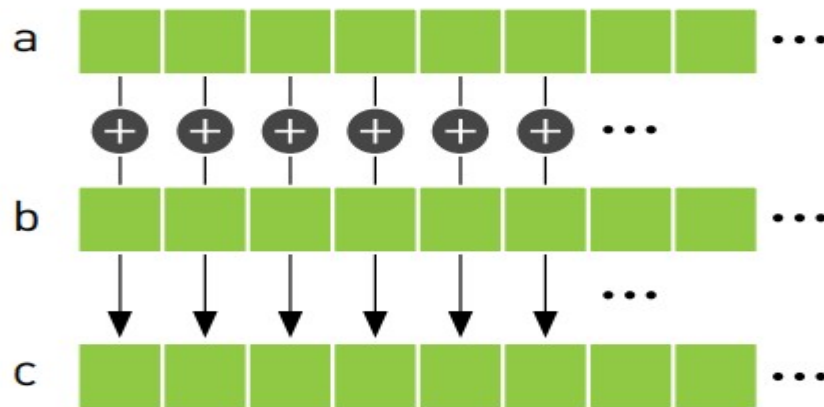


Figure 4.1 Summing two vectors

```

#include "../common/book.h"

#define N 10

void add( int *a, int *b, int *c ) {
    int tid = 0;    // this is CPU zero, so we start at zero
    while (tid < N) {
        c[tid] = a[tid] + b[tid];
        tid += 1;    // we have one CPU, so we increment by one
    }
}

int main( void ) {
    int a[N], b[N], c[N];

    // fill the arrays 'a' and 'b' on the CPU
    for (int i=0; i<N; i++) {
        a[i] = -i;
        b[i] = i * i;
    }

    add( a, b, c );
}

```

```
// display the results
for (int i=0; i<N; i++) {
    printf( "%d + %d = %d\n", a[i], b[i], c[i] );
}

return 0;
}
```

Most of this example bears almost no explanation, but we will briefly look at the `add()` function to explain why we overly complicated it.

```
void add( int *a, int *b, int *c ) {
    int tid = 0;    // this is CPU zero, so we start at zero
    while (tid < N) {
        c[tid] = a[tid] + b[tid];
        tid += 1;    // we have one CPU, so we increment by one
    }
}
```

```
void add( int *a, int *b, int *c ) {  
    for (i=0; i < N; i++) {  
        c[i] = a[i] + b[i];  
    }  
}
```

- one core initialize the loop with tid = 0 and another with tid = 1. The first core would add the even-indexed elements, and the second core would add the odd indexed elements.

CPU CORE 1

```
void add( int *a, int *b, int *c )
{
    int tid = 0;
    while (tid < N) {
        c[tid] = a[tid] + b[tid];
        tid += 2;
    }
}
```

CPU CORE 2

```
void add( int *a, int *b, int *c )
{
    int tid = 1;
    while (tid < N) {
        c[tid] = a[tid] + b[tid];
        tid += 2;
    }
}
```


GPU vector sums

```
#include "../common/book.h"

#define N 10

int main( void ) {
    int a[N], b[N], c[N];
    int *dev_a, *dev_b, *dev_c;

    // allocate the memory on the GPU
    HANDLE_ERROR( cudaMalloc( (void**)&dev_a, N * sizeof(int) ) );
    HANDLE_ERROR( cudaMalloc( (void**)&dev_b, N * sizeof(int) ) );
    HANDLE_ERROR( cudaMalloc( (void**)&dev_c, N * sizeof(int) ) );

    // fill the arrays 'a' and 'b' on the CPU
    for (int i=0; i<N; i++) {
        a[i] = -i;
        b[i] = i * i;
    }
}
```

```
// copy the arrays 'a' and 'b' to the GPU
HANDLE_ERROR( cudaMemcpy( dev_a, a, N * sizeof(int),
                          cudaMemcpyHostToDevice ) );

HANDLE_ERROR( cudaMemcpy( dev_b, b, N * sizeof(int),
                          cudaMemcpyHostToDevice ) );

add<<<N,1>>>( dev_a, dev_b, dev_c );

// copy the array 'c' back from the GPU to the CPU
HANDLE_ERROR( cudaMemcpy( c, dev_c, N * sizeof(int),
                          cudaMemcpyDeviceToHost ) );

// display the results
for (int i=0; i<N; i++) {
    printf( "%d + %d = %d\n", a[i], b[i], c[i] );
}

// free the memory allocated on the GPU
cudaFree( dev_a );
cudaFree( dev_b );
cudaFree( dev_c );

return 0;
}
```

You will notice some common patterns that we employ again:

- We allocate three arrays on the device using calls to `cudaMalloc()`: two arrays, `dev_a` and `dev_b`, to hold inputs, and one array, `dev_c`, to hold the result.
- Because we are environmentally conscientious coders, we clean up after ourselves with `cudaFree()`.
- Using `cudaMemcpy()`, we copy the input data to the device with the parameter `cudaMemcpyHostToDevice` and copy the result data back to the host with `cudaMemcpyDeviceToHost`.
- We execute the device code in `add()` from the host code in `main()` using the triple angle bracket syntax.

Moving on, our `add()` routine looks similar to its corresponding CPU implementation:

```
__global__ void add( int *a, int *b, int *c ) {  
    int tid = blockIdx.x;    // handle the data at this index  
    if (tid < N)  
        c[tid] = a[tid] + b[tid];  
}
```

```
__global__ void add( int *a, int *b, int *c ) {  
    int tid = blockIdx.x;    // handle the data at this index  
    if (tid < N)  
        c[tid] = a[tid] + b[tid];  
}
```

- N as the number of parallel blocks.
- Collection of parallel blocks a grid.
- This specifies to the runtime system that we want a one-dimensional grid of N blocks (scalar values are interpreted as one-dimensional).
- These threads will have varying values for `blockIdx.x`, the first taking value 0 and the last taking value N-1.
- Ex: imagine four blocks, all running through the same copy of the device code but having different values for the variable `blockIdx.x`.
- See the figures: Actual code being executed in each of the four parallel blocks looks like after the runtime substitutes the appropriate block index for `blockIdx.x`

BLOCK 1

```
__global__ void  
add( int *a, int *b, int *c ) {  
    int tid = 0;  
    if (tid < N)  
        c[tid] = a[tid] + b[tid];  
}
```

BLOCK 2

```
__global__ void  
add( int *a, int *b, int *c ) {  
    int tid = 1;  
    if (tid < N)  
        c[tid] = a[tid] + b[tid];  
}
```

BLOCK 3

```
__global__ void  
add( int *a, int *b, int *c ) {  
    int tid = 2;  
    if (tid < N)  
        c[tid] = a[tid] + b[tid];  
}
```

BLOCK 4

```
__global__ void  
add( int *a, int *b, int *c ) {  
    int tid = 3;  
    if (tid < N)  
        c[tid] = a[tid] + b[tid];  
}
```

- `HANDLE_ERROR()` macros that we've sprinkled so liberally throughout the code will detect and alert you to the situation.